

Using Description Logics for RDF Constraint Checking and Closed-World Recognition

Peter F. Patel-Schneider

Nuance Communications

1198 East Arques Avenue, Sunnyvale, California, U. S. A.

pfpschneider@gmail.com

Abstract

RDF and Description Logics work in an open-world setting where absence of information is not information about absence. Nevertheless, Description Logic axioms can be interpreted in a closed-world setting and in this setting they can be used for both constraint checking and closed-world recognition against information sources. When the information sources are expressed in well-behaved RDF or RDFS (i.e., RDF graphs interpreted in the RDF or RDFS semantics) this constraint checking and closed-world recognition is simple to describe. Further this constraint checking can be implemented as SPARQL querying and thus effectively performed.

There has recently been considerable attention paid to the problem of validating RDF (Cyganiak, Wood, and Lanthaler 2014) or RDFS information. There are several commercial systems that provide facilities for RDF validation, including TopQuadrant's SPIN (TopQuadrant 2011) and Clark&Parsia's Stardog ICV (Pérez-Urbina, Sirin, and Clark 2012; Clark&Parsia 2014). There are several proposals for specifying the desired form of RDF information, such as Resource Shapes (Ryman 2014). In 2013 W3C held an RDF Validation Workshop (W3C 2013) to gauge interest in the area, and W3C has started a new working group on RDF validation (W3C 2014).

Just what, however, is RDF validation?

Some accounts and systems (such as Stardog ICV) identify RDF validation with satisfying integrity constraints, similar to checking database integrity constraints. In this account, there are conditions (constraints) placed on instances of classes, such as requiring that every person has a name and an address, both strings. The defining characteristic here is that explicit information is needed to satisfy the integrity constraint. To pass the constraint that a person has a name it is necessary to provide a particular string for the name of the person, and not just state that the person has some unknown name.

Other accounts and proposals, such as OSLC Resource Shapes (Ryman, Hors, and Speicher 2013) (used by the Open Services for Lifecycle Collaboration community) and ShEx (Prud'hommeaux 2014; Solbrig and Prud'hommeaux

2014), identify RDF validation more with recognition, similar to determining whether an individual belongs to a Description Logic (Baader et al. 2010) description. For example one might define shape that requires a name and an address and then ask which individuals satisfy the shape's constraint. Here, in contrast to the previous situation, the validation is divorced from any type information in the data. Again, however, there is the requirement that explicit information is needed to match the shape—a particular name must be provided, not just information that there must be one.

As shown in this paper, Description Logics can be used to provide the necessary framework for both checking constraints and providing closed-world recognition facilities, and thus cover most of what SPIN and ShEx provide.

Why then are there claims (Ryman, Hors, and Speicher 2013; Fokoue and Ryman 2013) that OWL (Motik, Patel-Schneider, and Parsia 2012)—the Semantic Web Description Logic—is inadequate for these purposes? There are several aspects of the standard view of Description Logics that might not be consonant with constraints and the kind of recognition that might be desired. However, Description Logic syntax and semantics, and their instantiation in OWL, can serve as the basis for RDF constraint checking and closed-world recognition. The only change required is to consider a closed-world variation of the Description Logic semantics. Then the development of RDF constraint checking and closed-world recognition is easy.

The Basic Idea

Closed-World Recognition

Let's first look at recognition. In recognition we want to determine whether a particular node in an RDF graph matches some criteria. For example, John in the RDF graph¹

```
ex:John foaf:name "John"^^xsd:string.      (1)
ex:John foaf:phone "+19085551212"^^xsd:string
.
ex:John exo:friend ex:Bill .
ex:John exo:friend ex:Willy .
```

¹Turtle (Prud'hommeaux and Carothers 2014) will be used for writing RDF graphs throughout this paper. Prefix and base statements will generally be omitted.

matches the ShEx shape

$$\{ \text{foaf:name } \text{xsd:string}, \quad (2) \\ \text{foaf:phone } \text{xsd:string}, \\ \text{exo:friend } [2] \}$$

because John has a string value for his name, a string value for his phone, and two friends.

Determining whether an individual belongs to a Description Logic concept is also recognition. The Description Logic description² corresponding to the ShEx Shape (2) is

$$\begin{aligned} &= 1 \text{ foaf:name } \sqcap \forall \text{foaf:name.xsd:string } \sqcap \quad (3) \\ &= 1 \text{ foaf:phone } \sqcap \forall \text{foaf:phone.xsd:string } \sqcap \\ &= 2 \text{ exo:friend} \end{aligned}$$

So it seems that Description Logics can easily handle ShEx recognition. Although the ShEx syntax is somewhat more compact here, as ShEx has constructs that combine number restrictions and value restrictions, the Description Logic syntax is not verbose and is quite reasonable.

However, John does not match (3) in the standard semantics of Description Logic. This is precisely because in this standard reading, *and in RDF*, the absence of information is not information about absence. In the standard Description Logic reading, and also in RDF, John could have more than one name as far as the information in the above RDF graph is concerned. Many Description Logics (and, again, RDF too) also do not assume that different names refer to different individuals. So Bill and Willy could be the same person as far as the information in the above RDF graph is concerned.

It turns out that expressive Description Logics have facilities to explicitly state information about absence and information about differences and thus can be used to state complete information, at least on a local level. For example, if we add information to (1) stating that John has only one name and phone, that John's only friends are Bill and Willy, *and* that Bill is not the same as Willy, as in

$$\begin{aligned} \text{ex:John} &\in \leq 1 \text{ foaf:name} \\ \text{ex:John} &\in \leq 1 \text{ foaf:phone} \\ \text{ex:John} &\in \forall \text{exo:friend}.\{\text{ex:Bill}, \text{ex:Willy}\} \\ \text{ex:Bill} &\neq \text{ex:Willy} \end{aligned}$$

then John does match (3).

So it is not that Description Logics (including OWL) do not perform recognition as in ShEx, it is just that Description Logics do not make the assumption that absence of information is information about absence. In expressive Description Logics (again including OWL) it is possible to explicitly state what comes implicitly from the assumption that absence of information is information about absence.

However, suppose that we want to make this assumption generally? We could manually add a lot of axioms like the ones above, but this is both tedious and error-prone, and thus not at all a viable solution. Instead we can proceed by making the assumption that if the truth of some fact cannot be determined from the information given, then that fact is false.

²The abstract syntax (Baader et al. 2010) for Description Logics—a compact but non-ASCII syntax—will be used throughout this paper.

This is often called the closed world assumption or negation by failure, as the failure to prove some fact is used to support its falsity. There is a very large body of work on this topic (see the Related Work section for pointers into this work) and there are many tricky questions that arise with respect to closure in any sophisticated formalism, and expressive Description Logics (including OWL) are indeed sophisticated. As well, reasoning in expressive Description Logics that also have closed world facilities is extremely difficult, even in simple cases.

Fortunately RDF and RDFS are unsophisticated and inexpressive, so neither the tricky questions nor the reasoning difficulties arise if all information comes in the form of RDF triples interpreted under the RDF or RDFS semantics (Hayes and Patel-Schneider 2014). The basic idea is to treat the triples (and their RDF or RDFS consequences, if desired) as completely describing the world. In this treatment

1. if a triple is not present then it is false and
2. different IRIs denote different individuals.

This is precisely the same idea that underlies model checking, where a model is a finite set of ground first-order facts and everything else is false. First-order inference is undecidable, but determining whether a first-order sentence is true in one particular model (model checking) is much, much easier.

In this way it is possible to use the Description Logic syntactic and semantic machinery to define how to recognize descriptions under the same assumptions that underlie ShEx. The only change from the standard Description Logic setup is to define how to go from an RDF graph to the Description Logic model that the RDF graph embodies. Definitions, even recursive definitions, can also be handled.

This is all quite easy and conforms to a common thread of both theoretical and practical work. It also matches the theoretical underpinning of Stardog ICV (Pérez-Urbina, Sirin, and Clark 2012; Clark&Parsia 2014). Further, the approach can be implemented by translation into SPARQL queries, showing that it is practical. (There may be some constructs of very expressive description logics that do not translate into SPARQL queries when working with complete information, but at least the Description Logic constructs that correspond to the usual recognition conditions do so translate.)

Constraint Checking

Constraint checking does not appear to be part of the services provided by Description Logics. This has lead to claims that OWL cannot be used for constraint checking. However inference, which is the core service provided by Description Logics, and constraint checking are indeed very closely related.

Inference is the process of determining what follows from what has been stated. Inference ranges from simple (John is a student, students are people, therefore John is a person) to the very complex (involving reasoning by cases, *reductio ad absurdum*, or even noticing that an infinite sequence of inferences will not produce any useful information). Determining whether a constraint holds is just determining whether the constraint follows from the given information.

Figure 1: Data for Example

<code>ex:Amy rdf:type exo:UniStudent .</code> <code>ex:Amy foaf:name "Amy"^^xsd:string .</code> <code>ex:Bill rdf:type exo:UniStudent .</code> <code>ex:Bill foaf:name "Bill"^^xsd:string .</code> <code>ex:John rdf:type exo:GrStudent .</code> <code>ex:John foaf:name "John"^^xsd:string .</code> <code>ex:John exo:supervisor ex:Len .</code> <code>ex:Susan rdf:type exo:Person .</code> <code>ex:Susan foaf:name "Susan"^^xsd:string .</code> <code>ex:Len rdf:type exo:Faculty .</code> <code>ex:Len foaf:name "Len"^^xsd:string .</code>	<code>ex:Amy exo:enrolled ex:SUNYOrange .</code> <code>ex:Bill exo:enrolled ex:ReindeerPoly .</code> <code>ex:Bill exo:enrolled ex:HudsonValley .</code> <code>ex:John exo:enrolled ex:ReindeerPoly .</code> <code>ex:Susan exo:enrolled ex:ReindeerPoly .</code> <code>ex:Susan exo:enrolled ex:SUNYOrange .</code> <code>ex:Susan exo:enrolled ex:HudsonValley .</code> <code>ex:Len exo:affiliation ex:ReindeerPoly .</code> <code>ex:Len exo:affiliation ex:SUNYOrange .</code> <code>ex:SUNYOrange rdf:type exo:ResOrg .</code> <code>ex:HudsonValley rdf:type exo:Uni .</code>	<code>ex:Amy exo:friend ex:Bill .</code> <code>ex:Amy exo:friend ex:John .</code> <code>ex:Bill exo:friend ex:Amy .</code> <code>ex:Bill exo:friend ex:John .</code> <code>ex:John exo:friend ex:Amy .</code> <code>ex:John exo:friend ex:Bill .</code> <code>ex:John exo:friend ex:Len .</code> <code>ex:Len exo:friend ex:Amy .</code> <code>ex:Len exo:friend ex:Susan .</code>
--	--	--

Figure 2: RDFS Ontology for Example

<code>foaf:name rdfs:range xsd:string .</code> <code>exo:UniStudent rdfs:subClassOf exo:Person .</code> <code>exo:GrStudent rdfs:subClassOf exo:UniStudent .</code> <code>exo:Faculty rdfs:subClassOf exo:Person .</code> <code>exo:Uni rdfs:subClassOf exo:Organization .</code> <code>exo:ResOrg rdfs:subClassOf exo:Organization .</code>	<code>exo:enrolled rdfs:domain exo:UniStudent .</code> <code>exo:enrolled rdfs:range exo:Uni .</code> <code>exo:supervisor rdfs:domain exo:GrStudent .</code> <code>exo:supervisor rdfs:range exo:Faculty .</code> <code>exo:affiliation rdfs:domain exo:Person .</code> <code>exo:affiliation rdfs:range exo:Organization .</code>
---	--

Figure 3: Constraints and Recognition Axioms for Example

1 <code>exo:Person \sqcap exo:Organization $\equiv \{\}$</code> 2 <code>exo:Person \sqsubseteq = 1 foaf:Name $\sqcap \forall$ foaf:Name.xsd:string</code> 3 <code>exo:UniStudent \sqsubseteq \geq 1 exo:enrolled $\sqcap \forall$ exo:enrolled.exo:Uni</code> 4 <code>exo:GrStudent \sqsubseteq = 1 exo:enrolled $\sqcap \forall$ exo:enrolled.exo:ResOrg</code> 5 <code>exo:Faculty \sqsubseteq \geq 1 exo:affiliation $\sqcap \forall$ exo:affiliation.exo:Uni</code> 6 <code>exo:Faculty \sqsubseteq \leq 1 exo:affiliation.exo:ResOrg</code>	7 <code>exo:Faculty \sqsubseteq \leq 5 exo:supervisor . exo:GrStudent</code> 8 <code>exo:Uni \sqsubseteq \geq 2 exo:enrolled</code> 9 <code>exo:GrStudent exo:enrolled \sqsubseteq</code> <code> exo:supervisor \circ exo:affiliation</code> <code>ex:HecticStudent \equiv \geq 3 exo:enrolled</code> <code>ex:StudentFriend \equiv \geq 2 exo:friend . ex:StudentFriend</code>
---	---

Again, however, constraint checking is generally done with respect to complete information. So, to determine whether the constraint

$$\begin{aligned}
& ex:John \in exo:Person \sqcap \\
& \quad = 1 \text{ foaf:name } \sqcap \forall \text{ foaf:name.xsd:string } \sqcap \\
& \quad = 1 \text{ foaf:phone } \sqcap \forall \text{ foaf:phone.xsd:string } \sqcap \\
& \quad = 2 \text{ exo:friend}
\end{aligned} \tag{4}$$

is valid in the presence of (locally) complete information such as

```

ex:John  $\in$  exo:Person
ex:John foaf:name "John"^^xsd:string .
ex:John  $\in$  = 1 foaf:name
ex:John foaf:phone "+19085551212"^^xsd:string .
ex:John  $\in$  = 1 foaf:phone
ex:John exo:friend ex:Bill .
ex:John exo:friend ex:Willy .
ex:John exo:friend ex:Susan .
ex:John  $\in$   $\forall$  exo:friend. { ex:Bill, ex:Willy, ex:Susan }
ex:Bill  $\neq$  ex:Willy
ex:Bill  $\neq$  ex:Susan
ex:Susan  $\neq$  ex:Willy

```

is simply a matter of determining whether the constraint follows from the information. (Generally constraints like (4) are written to handle all the members of a class as in

$$\begin{aligned}
& exo:Person \sqsubseteq \\
& \quad = 1 \text{ foaf:name } \sqcap \forall \text{ foaf:name.xsd:string } \sqcap \\
& \quad = 1 \text{ foaf:phone } \sqcap \forall \text{ foaf:phone.xsd:string } \sqcap \\
& \quad = 2 \text{ exo:friend}
\end{aligned}$$

instead of just a single node, but the principle is the same.)

So a way to do constraints in Description Logics is to first set up complete information, and then just perform inference. This approach has been explored in the context of letting certain roles be completely specified as in a database (Patel-Schneider and Franconi 2012). Other approaches to constraints in Description Logics (de Bruijn et al. 2005; Motik, Horrocks, and Sattler 2009; Tao et al. 2010; Donini, Nardi, and Rosati 2002; Sengupta, Krisnadhi, and Hitzler 2011) are considerably more complex, as they deal with the complexities that arise when there are multiple ways to complete the information. However, all these approaches largely agree when there is only a single way to complete the information.

Setting up complete information is just what was done above for closed-world recognition, so this technique can also be used for constraint checking. Of course, this doesn't mean that you have to implement Description Logic inference with complete information the same way that you need to with incomplete information. In fact, as above, constraint checking can be implemented as SPARQL queries.

Example

Here is a small example of how Description Logic constructs can be used for constraint checking. There are three separate kinds of information in the example. The RDF triples in Figure 1 provide the data for the example. The RDFS ontology in Figure 2 provides the organization of the data. The Description Logic axioms in Figure 3 provide the constraints to be validated against the data and the ontology and the classes vocabulary for closed-world recognition.

The constraints are all satisfied, except for one, as follows:

1. Nothing can be inferred to be both a person and an organization, and so persons and organizations are disjoint.
2. Every object that can be inferred to be a person (i.e., students and faculty) has a single name provided, and that name is a string, so every person has exactly one name in the closure.
3. All students (and grad students) are enrolled in universities—the range of *exo:enrolled* is *exo:Uni*, which makes the typing part of the constraint redundant here.
4. Reindeer Poly is not specified to be a research organization, so although John is enrolled exactly once the constraint on graduate students being enrolled in research organizations is not satisfied.
5. All faculty (Len) are affiliated with only universities.
6. All faculty (again only Len) are affiliated with at most one research organization, as Reindeer Poly is not specified to be a research organization and Len is only affiliated with SUNY Orange and Reindeer Poly.
7. All faculty supervise fewer than five graduate students, as the only faculty (Len) only supervises one student (John).
8. Each university (SUNY Orange, Reindeer Poly, and Hudson Valley) has at least two students enrolled in it, because Amy, Bill, John, and Susan are different individuals.
9. For every graduate student enrollment (*exo:enrolled* domain-restricted to *exo:GrStudent*) there is a supervisor of the graduate student affiliated with the university. This constraint uses an auxiliary non-recursive equivalence definition.

The only hectic student is Susan, as she is the only person with at least three enrollments. However, Amy, Bill, and John all belong to *ex:StudentFriend* because when maximally interpreting *ex:StudentFriend* they each have at least two friends who belong to *ex:StudentFriend*. Len does not belong to *ex:StudentFriend* even though he has two friends, because Susan has too few friends and cannot belong to *ex:StudentFriend*.

One might want to validate that domain and range types are not inferred, but are instead explicitly stated in the data. This can be done by using a version of the ontology without the domain and range statements and validating against a set of constraints that just have the removed domain and range statements. In the example, this would detect that Susan was not stated to be a student, violating the domain constraint for *exo:enrolled*; that SUNY Orange and Reindeer Poly were not stated to be universities, violating the

range constraint for *exo:enrolled*; and that Reindeer Poly was not stated to be an organization, violating the range constraint for *exo:affiliation*. All other domain and range constraints would be satisfied, as some required class memberships would be inferred from the subclass statements.

Related Work

The closest work in a technical sense is the work of Patel-Schneider and Franconi (2012). In that work some properties and classes were considered as closed, which turned description logic axioms involving those properties and classes into constraints. RDF and RDFS are very similar to a situation where all properties and classes are closed. The current paper adds the idea of closed-world recognition, which was only implicit in the previous work, and maximal extensions, which provide a much better treatment of recursive definitions, particularly in the monotone case.

The work of Motik, Horrocks, and Sattler (2009) and of Tao *et al.* (2010) both divide up axioms into regular axioms and constraints. They both also permit general Description Logic axioms, not just RDF or RDFS graphs as here. To handle full Description Logic information requires a much more complex construction, involving minimal interpretations. Neither consider closed-world recognition. Tao *et al.* use SPARQL queries as a partial translation of their constraints and forms a basis for Stardog ICV.

The work of Sengupta, Krisnadhi, and Hitzler (2011) uses circumscription as the mechanism to minimize interpretations. It is otherwise similar to the previous efforts. The work of Donini, Nardi, and Rosati (2002) uses autoepistemic constructs within axioms to model constraints, and is thus quite different from the approach here. OWL Flight (de Bruijn *et al.* 2005) is a subset of OWL where axioms are given meaning as Datalog constraints. Again, as an expressive Description Logic is handled the construction is more complex than the one here. RDFUnit (Kontokostas *et al.* 2014) has a component that turns RDFS axioms and simple OWL axioms into SPARQL queries that check for data that does not match the axiom and so is somewhat similar to this work. However, there is no notion that RDFUnit is turning ontology axioms into constraints that cover the entire meaning of the axiom.

Shex (Solbrig and Prud'hommeaux 2014) uses very different mechanisms. It builds up shapes, which are akin to definitions of classes, and gives them meaning by a translation into a recursive extension of Z over an abstraction of RDF graphs. Entire documents or document portions are then matched against these shapes.

The Details, but Not All the Details

Description Logic Semantics

The semantics of Description Logics are generally given as a model theory, as for OWL DL (Motik, Patel-Schneider, and Parsia 2012). OWL DL has a complex semantics, as far as Description Logics go, to cover all its constructs and to make it more compatible with RDF. The semantics here will follow the semantics of OWL, with the exception that any property can have both individuals, e.g., *ex:John*, and data values, e.g., "John"^^xsd:string, as values.

The fundamental building block of Description Logics semantics is interpretations, which provide a meaning for the primitive constructs in terms of a particular domain of discourse. The meaning of an individual name, such as *ex:John*, is an element of this domain, here that element that we think of as being John. Literal values, such as *"John"^^xsd:string*, are treated specially—their meaning is determined by their datatype. The meaning (or interpretation) of a named concept such as *exo:person*, is a subset of the domain, here those individuals that we might think of as people. The meaning of a named property, such as *exo:friend*, is a set of pairs over the domain, here those pairs that we might think of as being the friend-of relationship. The meaning of non-primitive constructs, such as the description \exists *exo:friend*, are built up from these primitives, resulting here in the set of domain elements that are related to exactly two domain elements via the meaning of *exo:friend*.

Axioms, such as *ex:John* \in *exo:Person*, are true precisely when the meaning of their parts satisfies a particular relationship, here that the meaning of *ex:John* is an element of the meaning of *exo:Person*. There are some other aspects to this simple story, to handle the differences between individuals, e.g., *ex:John*, and data values, e.g., *"John"^^xsd:string*, and to make reasoning over some constructs easier. A Description Logic model of a set of axioms (including what we might call facts), is then just an interpretation that makes all the axioms true.

Canonical Interpretations of RDF Graphs

In an interpretation everything is specified, so each interpretation has complete information. The basic idea is thus to construct an interpretation making just the triples in an RDF graph true and then work with that single interpretation. In this way information that is absent from the RDF graph is considered false.

We can think of most RDF graphs as sets of Description Logic axioms, particularly as we are ignoring the common Description Logic division of properties into properties that have objects that are individuals and properties that have objects that are data values.³ This correspondence breaks down in two areas: 1/ when the built-in RDF and RDFS vocabulary is used in unusual ways (e.g., making *rdfs:subClassOf* a sub-property of *rdfs:subPropertyOf*), and 2/ if reasoning about individuals can affect reasoning about classes (e.g., forcing two individuals that are also classes to be the same). The abuse and extension of the built-in RDF and RDFS vocabulary is rare, so we just exclude these RDF graphs from our account. (Particular extensions of the RDF and RDFS vocabulary could be built in to an extension of the approach given here.) In RDF and RDFS, but not in OWL, reasoning about individuals can affect reasoning about classes. This is also rare, so we will not handle these inferences.

³This division has been made so that reasoners do not have to worry about data values having properties, but if we are constructing models this is not a problem. It is easy to revise the treatment here to bring back this division.

Definition 1 *Given an RDF graph G with no ill-formed literals and no triples stating membership in a datatype, we construct the canonical Description Logic interpretation of G as follows.*

1. *Datatypes are formed for all the datatypes in the graph, and given meaning in the usual way.*
2. *The domain of the interpretation consists of the non-literal nodes of the RDF graph plus the properties in the graph and the mapping for nodes is the identity mapping. (One might think that this is not an appropriate way to construct an interpretation, as it sets the meaning of *ex:John* to *ex:John*, not anything that we might think of as being John, but as far as the formal machinery is concerned, the actual domain elements are not important.)*
3. *The set of literal values is constructed in the usual way from the datatypes. An extra copy of the integers is added to ensure an infinite number of literal values.*
4. *Classes are formed for each node in the graph that has an *rdf:type* link with it as an object or belongs to *rdfs:Class* and their extensions are the set of nodes for which *rdf:type* triples link them to the class.*
5. *Properties (note that we are ignoring the Description Logic division of properties) are formed for each predicate in the graph and also for each node that belongs to *rdf:Property*, and their extensions are the set of pairs taken from triples in the graph with the property as predicate.*

So from the initial RDF graph in this paper, we end up with a canonical interpretation with six domain elements, *ex:John*, *ex:Bill*, *ex:Willy*, *foaf:name*, *foaf:phone*, and *exo:friend*. The interpretation of *foaf:name* consists of just $\langle \text{ex:John}, \text{"John"} \rangle$. The interpretation of *foaf:phone* consists of just $\langle \text{ex:John}, \text{"+19085551212"} \rangle$. The interpretation of *exo:friend* consists of $\langle \text{ex:John}, \text{ex:Bill} \rangle$ and $\langle \text{ex:John}, \text{ex:Willy} \rangle$.

For constraints and descriptions that use only vocabulary in the RDF graph all we do is work with this interpretation and consider whether the constraint axiom is true in this interpretation so the development is easy. It is obvious that *ex:John* belongs to the interpretation of the first Description Logic description given above, as expected.

Evaluating constraints on the canonical interpretation of a graph is essentially the same as evaluating them on the graph itself. Systems that evaluate constraints on an RDF graph, like ShEx, thus work in a manner very similar to the approach taken here.

Extending to New Classes

For closed-world recognition, it is useful to define new classes, as in

$$\begin{aligned} \text{ex:PurePerson} &\equiv \geq 1 \text{ exo:friend} \sqcap \\ &\quad \forall \text{exo:friend. ex:PurePerson} \end{aligned}$$

There are several possibilities for the meaning of new classes that are recursively defined. The new classes could be interpreted as broadly as possible, as narrowly as possible, or in any consistent manner.

It appears in ShEx that such classes as to be interpreted as broadly as possible. For example, in the RDF graph

ex:John *exo:friend* *ex:Bill* .
ex:Bill *exo:friend* *ex:John* .

the ShEx approach would be that both *ex:John* and *ex:Bill* belong to *ex:PurePerson*. We will take this approach here and interpret new classes as broadly as possible

New classes are handled by considering extensions of the interpretations defined as above. An extension of an interpretation is a new interpretation 1/ with the same domain as the original interpretation, and 2/ that has the same meaning for all individuals, named classes, and named properties in the original interpretation. The extension is allowed to have new named classes, but not new named properties or new individuals. New individuals are not allowed because some Description Logic constructs are sensitive to the set of individuals. New properties are not allowed because they may increase the computational complexity of closed-world recognition.

An interpretation is an extended canonical model of an RDF graph with respect to a set of constraints if it is an extension of the canonical model of the RDF graph and is a model of the constraints.

To interpret recursively defined classes as broadly as possible not all extended canonical models are considered, only maximal ones. A model is maximal among a set of models if there is no other model in the set that 1/ interprets all classes as supersets of their interpretation in the maximal model, and 2/ interprets at least one class as a strict superset of its interpretation in the maximal model.

An individual is recognized as belonging to a description if its interpretation belongs to the interpretation of the description in all maximal extended canonical models.

It turns out that there is only one (up to isomorphism) maximal extended canonical model of the above definition of *ex:PurePerson*. In this model both *ex:John* and *ex:Bill* are in the extension of *ex:PurePerson*. If all new classes are monotone in all the other new classes (i.e., if the extension of some class grows then no other class extensions shrink) then there is always exactly one maximal extension.⁴

RDF and RDFS Semantics

So everything looks fine. We go from an RDF graph to a slightly modified Description Logic interpretation and from there perform constraint checking by determining whether a set of Description Logic axioms are satisfied in the interpretation or in a set of maximal extensions of the interpretation. We can also perform closed-world recognition by determining the interpretation of the new defined classes in the axioms and these classes can even be defined recursively.

However, there is one missing part of the story. If the RDF graph includes triples that trigger RDF or RDFS inferences that are not already in the RDF graph the interpretation will not look like an RDF (or RDFS) interpretation. For example, if the RDF graph is

ex:John *rdf:type* *exo:Student* .
exo:Student *rdfs:subClassOf* *exo:Person* .

⁴Proof sketches of claims here and later in the paper are in the appendix of this extended version.

our canonical interpretation states that *ex:John* does not belong to *exo:Person*, which goes against the RDFS meaning of the above graph.

Fortunately, it is relatively easy to recover from this problem. All that is needed is to add all the RDF (or RDFS) consequences to the graph. Yes, there are an infinite number of these consequences, but our formal development does not care whether the graph is finite or infinite. For complexity analysis and implementation it is not hard to come up with a finite representation of these consequences, like the one initially done by ter Horst (2005), and make the minor fix-ups needed to determine correct answers from the answers gleaned from this finite approximation.

This all works because the RDF (or RDFS) consequences of an RDF graph can be represented as an RDF graph. OWL consequences cannot be so represented, as the consequences in OWL can be disjunctive. This requires working with minimal equality and minimal models or some other way to single out only the desired interpretations as in previous work on Description Logic constraints, making the formal development much harder and presenting many more choices that have to be justified.

Complexity

It is easy to see that checking axioms against an interpretation is polynomial, as long as there is no new vocabulary in the axioms or no recursive definitions. The formulae corresponding to the axioms are just model checked against the interpretation.

If there are monotone recursive definitions then checking constraints and performing closed-world recognition can be done using techniques from Datalog, such as magic sets. For example, the extension of a recursively-defined class can first be computed ignoring the recursive portion. Violations of the recursive portion can then be checked and objects iteratively removed from the class.

These techniques cannot be used for non-monotone recursive definitions, as expanding one class or property might reduce another.

Implementation

The work of Tao *et al.* (2010) shows that the standard Description Logic constraints can be partly implemented as SPARQL queries when no new vocabulary is used. Tao *et al.* worked in a general OWL setting, where their approach is sound but not complete, but in an RDF setting the approach is both sound and complete, because there is only a single model that needs to be considered. This approach forms the basis of Stardog ICV (Clark&Parsia 2014). Indeed Stardog ICV is an implementation of the approach described in this paper showing how the approach to constraints here can be implemented by a translation to SPARQL. The work here can thus also be thought of as a simpler definition of the underpinning of Stardog ICV. Recent work at Mannheim by Thomas Bosch (see <https://github.com/boschthomas/OWL2-SPIN-Mapping>) translates OWL descriptions interpreted as constraints into

SPARQL using a similar approach, providing a different implementation.

Non-recursive closed-world recognition can be handled by using nested or repeated SPARQL queries. Monotone recursive closed-world recognition can be implemented using Datalog techniques. Non-monotone recursive closed-world recognition is more complex and cannot be handled in the same way. This indicates that excluding non-monotone recursive closed-world recognition could be a reasonable stance to take.

Conclusion

Description Logics can indeed be used for both the syntax and semantics of constraint checking and closed-world recognition in RDF, by employing an analogue of model checking, and much of both constraint checking and closed-world recognition can be effectively implemented using a translation to SPARQL queries. The main difference between closed-world recognition and constraint checking is that the former either has no axioms or only uses axioms defining names that do not occur in the RDF graph whereas constraint checking uses axioms that relate concepts appearing in the RDF graph to descriptions.

By restricting our information to be RDF or RDFS, i.e., working in situations where there is a unique minimal model, we obtain a simpler formulation, easy implementation, and good performance as compared to previous work in this area. The approach here can be easily extended to other subsets of OWL that have a unique minimal model. An OWL profile with this property is OWL RL (Motik et al. 2012).

References

- [Baader et al. 2010] Baader, F.; Calvanese, D.; McGuinness, D. L.; Nardi, D.; and Patel-Schneider, P. F., eds. 2010. *The Description Logic Handbook: Theory, implementation, and applications*. Cambridge University Press, 2nd edition.
- [Clark&Parsia 2014] Clark&Parsia. 2014. Stardog integrity constraint validation. <http://docs.stardog.com/icv>.
- [Cyganiak, Wood, and Lanthaler 2014] Cyganiak, R.; Wood, D.; and Lanthaler, M. 2014. RDF 1.1 concepts and abstract syntax. W3C Recommendation, <http://www.w3.org/TR/rdf-concepts/>.
- [de Bruijn et al. 2005] de Bruijn, J.; Polleres, A.; Lara, R.; and Fensel, D. 2005. OWL DL vs. OWL Flight: Conceptual modeling and reasoning for the semantic web. In *Proceedings of the Fourteenth World Wide Web Conference*, 623–632. ACM Press.
- [Donini, Nardi, and Rosati 2002] Donini, F. M.; Nardi, D.; and Rosati, R. 2002. Description logics of minimal knowledge and negation as failure. *ACM Transactions on Computational Logic* 3(2):177–225.
- [Fokoue and Ryman 2013] Fokoue, A., and Ryman, A. 2013. OSLC resource shape: A linked data constraint language. In W3C (2013). http://www.w3.org/2001/sw/wiki/images/b/b7/\RDFVal_Fokoue_Ryman.pdf.
- [2014] Hayes, P., and Patel-Schneider, P. F. 2014. RDF 1.1 semantics. W3C Recommendation, <http://www.w3.org/TR/rdf11-mt/>.
- [2014] Kontokostas, D.; Westphal, P.; Auer, S.; Hellmann, S.; Lehmann, J.; Cornelissen, R.; and Zaveri, A. J. 2014. Test-driven evaluation of linked data quality. In *Proceedings of the Twenty-Third World Wide Web Conference*, 747–758. ACM Press.
- [2012] Motik, B.; Grau, B. C.; Horrocks, I.; Wu, Z.; Fokoue, A.; and Lutz, C. 2012. OWL 2 Web Ontology Language: Profiles (second edition). W3C Recommendation, <http://www.w3.org/TR/owl2-profiles>.
- [2009] Motik, B.; Horrocks, I.; and Sattler, U. 2009. Bridging the gap between OWL and relational databases. *Journal of Web Semantics* 7(2):74–119.
- [2012] Motik, B.; Patel-Schneider, P. F.; and Parsia, B. 2012. OWL 2 web ontology language: Structural specification and functional-style syntax. W3C Recommendation, <http://www.w3.org/TR/owl2-syntax/>.
- [2012] Patel-Schneider, P. F., and Franconi, E. 2012. Ontology constraints in incomplete and complete data. In *Proceedings of the Eleventh International Semantic Web Conference*, 444–459. Springer.
- [2012] Pérez-Urbina, H.; Sirin, E.; and Clark, K. 2012. Validating RDF with OWL integrity constraints. <http://docs.stardog.com/icv/ics-specification.html>.
- [2014] Prud’hommeaux, E., and Carothers, G. 2014. RDF 1.1 turtle: Terse RDF triple language. W3C Recommendation, <http://www.w3.org/TR/turtle/>.
- [2014] Prud’hommeaux, E. 2014. Shape expressions 1.0 primer. W3C Member Submission, <http://www.w3.org/Submission/shex-primer/>.
- [2013] Ryman, A. G.; Hors, A. J. L.; and Speicher, S. 2013. OSLC resource shape: A language for defining constraints on linked data. In Bizer, C.; Heath, T.; Berners-Lee, T.; Hausenblas, M.; and Auer, S., eds., *Proceedings of the WWW2013 Workshop on Linked Data on the Web*. <http://ceur-ws.org/Vol-996/papers/ldow2013-paper-02.pdf>.
- [2014] Ryman, A. 2014. Resource shape 2.0. W3C Member Submission, <http://www.w3.org/Submission/shapes/>.
- [2011] Sengupta, K.; Krisnadhi, A.; and Hitzler, P. 2011. Local closed world semantics: Grounded circumscription for OWL. In *Proceedings of the Tenth International Semantic Web Conference*, 617–632. Springer.
- [2014] Solbrig, H., and Prud’hommeaux, E. 2014. Shape expressions 1.0 definition. W3C Member Submission, <http://www.w3.org/Submission/shex-defn/>.
- [2010] Tao, J.; Sirin, E.; Bao, J.; and McGuinness, D. L. 2010. Integrity constraints in OWL. In *Proceedings of the Twenty-Fourth National Conference on Artificial Intel-*

ligence, 1443–1448. American Association for Artificial Intelligence.

[2005] ter Horst, H. J. 2005. Completeness, decidability and complexity of entailment for RDF Schema and a semantic extension involving the OWL vocabulary. *Journal of Web Semantics* 3(2-3):79–115.

[2011] TopQuadrant. 2011. SPIN—modeling vocabulary. W3C Member Submission, <http://www.w3.org/Submission/spin-modeling/>.

[2013] W3C. 2013. W3C RDF validation workshop report. <http://www.w3.org/2012/12/rdf-val/>.

[2014] W3C. 2014. W3C RDF data shapes charter. <http://www.w3.org/2014/data-shapes/charter>.

Appendix

Description Logic Semantics

This is a condensed version of Description Logic semantics, largely taken from the OWL 2 semantics (Motik, Patel-Schneider, and Parsia 2012). Some parts of the semantics have been removed (including facets and naming) so that this account is easier to read. As well, the division of properties between object properties and data properties has been removed so that RDF properties that have both objects and data values as objects can be handled. Many notions from RDF and Description Logics will be used without definition.

Definition 2 A datatype consists of the name of the datatype, which is an IRI; the set of values for the datatype, and a partial mapping from strings to these values.

Definition 3 A (class, property, and individual) vocabulary is a tuple $\langle V_C, V_P, V_I \rangle$ where V_C , the classes, V_P , the properties, and V_I , the individuals, are each sets of IRIs and blank nodes,

Note that there is no requirement that the classes, properties, and individuals in a vocabulary be pairwise disjoint.

Definition 4 Given a vocabulary V and set of datatypes D , an interpretation is a tuple $\langle \Delta_I, \Delta_D, \cdot^C, \cdot^P, \cdot^I \rangle$, where

- Δ_I is a non-empty set of objects, the domain of the interpretation,
- Δ_D is a set of data values, containing at least all the values for datatypes in D ,
- \cdot^C maps V_C into subsets of Δ_I ,
- \cdot^P maps V_P into subsets of $\Delta_I \times (\Delta_I \cup \Delta_D)$, and
- \cdot^I maps V_I into elements of Δ_I .

The semantics also uses \cdot^{DT} , which maps datatypes in D into the set of their values as specified by the datatype, and \cdot^{LT} , which maps literals into their values as specified by the datatypes. \cdot^C and \cdot^P are extended to class expressions and property expressions in the usual way.

Definition 5 A Description Logic axiom (i.e., an OWL axiom) is true in an interpretation in the usual way, with appropriate modifications made to eliminate the distinction between object and data properties. A model of a set of Description Logic axioms is an interpretation that makes all the axioms true.

Same-Vocabulary Constraints

Definition 6 The vocabulary for an RDF graph G is $V = \langle V_C, V_P, V_I \rangle$ where

$$\begin{aligned} V_C &= \{C \mid \exists s \langle s, \text{rdf:type}, C \rangle \in G\} \cup \\ &\quad \{C \mid \langle C, \text{rdf:type}, \text{rdfs:Class} \rangle \in G\}, \\ V_P &= \{P \mid \exists x, o \langle s, P, o \rangle \in G\} \cup \\ &\quad \{P \mid \langle P, \text{rdf:type}, \text{rdf:Property} \rangle \in G\}, \end{aligned}$$

and V_I is the set of non-literal nodes in G .

Definition 7 Given a set of datatypes D , the canonical interpretation for an RDF graph G using only datatypes in D is the interpretation $I = \langle \Delta_I, \Delta_D, \cdot^C, \cdot^P, \cdot^I \rangle$ over the vocabulary of G , $\langle V_C, V_P, V_I \rangle$, and the datatypes in D , where

- $\Delta_I = V_C \cup V_P \cup V_I$,
- Δ_D is the union of all the values for datatypes in D disjointly unioned with a copy of the integers,
- $c^C = \{s \mid \exists \langle s, \text{rdf:type}, c \rangle \in G\}$, for $c \in V_C$,
- $p^P = \{\langle s, o \rangle \mid \exists \langle s, p, o \rangle \in G\}$, for $p \in V_P$, and
- $i^I = i$, for $i \in V_I$.

The canonical interpretation of G includes all classes and properties in G as individuals. It also constructs Description Logic classes and properties for the built-in RDF and RDFS classes and properties such as *rdfs:Class*, *rdf:type*, and *rdfs:subClassOf*. This is not exactly what might be expected, does create a reasonable interpretation that matches RDF and RDFS intuitions closely.

Theorem 1 Given a set of datatypes D and an RDF graph G using only these datatypes, if G is closed under the RDF (RDFS) rules of inference then a minor adjustment to the canonical interpretation of G is a model of G under the RDF (RDFS) semantics.

Proof: Adjustments first have to be made to the canonical model to turn it into an actual RDF interpretation. The RDF domain is the union of the objects and data values of the canonical interpretation. Literals are placed into the datatypes they belong to.

The proof is then via a simple case-by-case analysis of each semantic condition on RDF (RDFS) models.

Definition 8 Given a set of datatypes D , an RDF graph G with vocabulary V , and a set of Description Logic axioms C over V acting as constraints, C is satisfied by G iff the canonical interpretation of G is a model for C .

Note that in many cases where C is satisfied by G , C will not follow from G . For example

$\{ \text{ex:Person} \sqsubseteq \leq 1 \text{ foaf:name} \}$

is satisfied by but does not follow from

$\text{ex:John} \text{ rdf:type } \text{ex:Person} .$
 $\text{ex:John} \text{ foaf:name } \text{"John"}^{\text{^^xsd:string}} .$

Theorem 2 Given an RDF graph G with vocabulary V and a set of constraints C over V and with no blank nodes, checking whether G satisfies C can be done in polynomial time for most Description Logics, including OWL.

Proof Sketch: For most axioms checking can be reduced to checking inclusion relationships between descriptions in the canonical model. Determining the extension of a description in a model involves checking some local conditions, which can be easily done in polynomial time because there are no choices to be made in the model, which specifies the extension of all named classes and properties. Some axioms are not reducible to checking inclusion relationships (e.g, key axioms) but are similarly easy to check in a model. The prohibition of blank nodes is to prevent the creation of sets of individual axioms that require checking for the presence of graph structures in the canonical model.

Blank nodes do not cause a problem here because the only blank nodes allowed in the constraints are blank nodes that also are in the graph, and these blank nodes are treated just the same as if they are IRIs. The price paid is that it is not possible to use blank nodes in the constraints to perform structure matching, for example to see if there is some individual that is related to another particular individual via two separate role chains.

Theorem 3 *Given an RDF graph G and vocabulary V and a set of constraints C over V determining whether G satisfies C can be done by a number of SPARQL queries on G polynomial in the size of C .*

Proof Sketch: Description Logic description or property expressions have obvious translations to SPARQL queries that when run on G produce the extension of the description. The translations are the ones used in Stardog ICV as described by Tao *et al.* (2010). The Description Logic axioms that check the relationship between two description or property expressions can be checked by creating a SPARQL query that is empty if and only if the axioms is satisfied by G , again as done by Tao *et al.* Other Description Logic axioms, such as transitivity and keys can be treated in similar ways.

Class-Extended Constraints

Definition 9 *Given a vocabulary $V = \langle V_C, V_P, V_I \rangle$ and a set of datatypes D , $I' = \langle \Delta'_I, \Delta'_D, \cdot^{C'}, \cdot^{P'}, \cdot^{I'} \rangle$, an interpretation over vocabulary $V' = \langle V'_C, V'_P, V'_I \rangle$, is an extension of an interpretation $I = \langle \Delta_I, \Delta_D, \cdot^C, \cdot^P, \cdot^I \rangle$ over V iff*

- $V_C \subseteq V'_C$, $V_P = V'_P$, $V_I = V'_I$,
- $\Delta_I = \Delta'_I$, $\Delta_D = \Delta'_D$,
- $c^C = c^{C'}$ for $c \in V_C$,
- $p^P = p^{P'}$ for $p \in V_P$, and
- $i^I = i^{I'}$ for $i \in V_I$.

Note that only the class vocabulary can be extended; the property and individual vocabularies remain the same. Also note that the domain is unchanged.

Theorem 4 *If I is a model of the RDF graph G with vocabulary V and I' is an extension of I , then I' is a model of G .*

Proof: I and I' only differ on the class vocabulary outside V , which does not affect whether an interpretation is a model of G even for Description Logic constructs that are sensitive to the individual vocabulary.

Definition 10 *Given a set of datatypes D , an RDF graph G with vocabulary V , and a set of Description Logic axioms C whose properties (individuals) are all properties (individuals) from V , C is satisfied by G iff there is an extension of the canonical interpretation of G that is a model for C .*

Note that each class in C that has an equality definition that does not directly or indirectly refer to itself has the same class extension in each extension of the canonical interpretation of G that is a model for C .

Definition 11 *Given a set of datatypes D , an RDF graph G with vocabulary V , and a set of Description Logic axioms C whose properties (individuals) are all properties (individuals) from V , o is in the closed-world class extension of c for o an individual in G and c a class in G or C iff $m^I(o) \in m^C(c)$ in each extension, m , of the canonical interpretation of G that is a model for C .*

Definition 12 *Given a vocabulary $V = \langle V_C, V_P, V_I \rangle$ and a set of datatypes D , $I' = \langle \Delta'_I, \Delta'_D, \cdot^{C'}, \cdot^{P'}, \cdot^{I'} \rangle$, an interpretation over V is bigger than or equal to an interpretation $I = \langle \Delta_I, \Delta_D, \cdot^C, \cdot^P, \cdot^I \rangle$ over V iff $c^C \subseteq c^{C'}$ for all $c \in V_C$, $p^P = p^{P'}$ for all $p \in V_P$, and $i^I = i^{I'}$ for all $i \in V_I$.*

Definition 13 *Given a set of datatypes D , an RDF graph G with vocabulary V , and a set of Description Logic axioms C whose properties (individuals) are all properties (individuals) from V , an interpretation $I = \langle \Delta_I, \Delta_D, \cdot^C, \cdot^P, \cdot^I \rangle$ over vocabulary V' is a maximal extension model of G and C iff*

- I is an extension of the canonical interpretation of G ,
- I is a model of C , and
- there is no model of G and C that is bigger than I .

Definition 14 *Given a set of datatypes D , an RDF graph G with vocabulary V , and a set of Description Logic axioms C whose properties (individuals) are all properties (individuals) from V , o is in the maximal closed-world class extension of c for o and c nodes in G or C iff $m^I(o) \in m^C(c)$ in each maximal extension model, m , of G and C .*

Definition 15 *Given a set of datatypes D , an RDF graph G with vocabulary V , a set of constraints C whose properties (individuals) are all properties (individuals) from V , and two class names in C but not in V , C_1 and C_2 , C_1 is monotone with respect to C_2 if whenever I_1 and I_2 are models of G that satisfy C and have the same extension for all classes in C except for C_1 and C_2 then if $I_1^C(C_2) \subseteq I_2^C(C_2)$ then $I_1^C(C_1) \subseteq I_2^C(C_1)$.*

Definition 16 *Given a set of datatypes D , an RDF graph G with vocabulary V , and a set of constraints C whose properties (individuals) are all properties (individuals) from V , then C is monotone iff all class names in C but not in V are monotone with respect to all class names in C but not in V .*

Theorem 5 *Given a set of datatypes D , an RDF graph G with vocabulary V , and a set of monotone constraints C over G and V , there is one maximal extension for G and C up to isomorphism.*

Proof: Take a model of G and C that has the biggest extension for some class in C but not in V . This model must also have maximal extensions for all the other classes in C but not in V because otherwise C would not be monotone. This model is thus, up to isomorphism, maximal.